

Robust Order Statistics based Ensembles for Distributed Data Mining

Kagan Tumer

NASA Ames Research Center
MS 269-2, Moffett Field, CA, 94035-1000
kagan@ptolemy.arc.nasa.gov

Joydeep Ghosh

Department of Electrical and Computer Engineering,
University of Texas, Austin, TX 78712-1084
ghosh@ece.utexas.edu

Abstract

Integrating the outputs of multiple classifiers via combiners or meta-learners has led to substantial improvements in several difficult pattern recognition problems. In the typical setting investigated till now, each classifier is trained on data taken or resampled from a common data set, or randomly selected partitions thereof, and thus experiences similar quality of training data. However, in distributed data mining involving heterogeneous databases, the nature, quality and quantity of data available to each site/classifier may vary substantially, leading to large discrepancies in their performance. In this chapter we introduce and investigate a family of meta-classifiers based on order statistics, for robust handling of such cases. Based on a mathematical modeling of how the decision boundaries are affected by order statistic combiners, we derive expressions for the reductions in error expected when such combiners are used. We show analytically that the selection of the median, the maximum and in general, the i^{th} order statistic improves classification performance. Furthermore, we introduce the trim and spread combiners, both based on linear combinations of the ordered classifier outputs, and empirically show that they are significantly superior in the presence of outliers or uneven classifier performance. So they can be fruitfully applied to several heterogeneous distributed data mining situations, specially when it is not practical or feasible to pool all the data in a common data warehouse before attempting to analyze it.

1 Mining of Distributed Data Sources

An implicit assumption in traditional statistical pattern recognition and machine learning algorithms is that the data to be used for model development is available as a single flat file. This assumption is valid for virtually all popular benchmark datasets such as those available from ELENA, Statlog or the UCI machine learning repository. Such datasets are small or medium sized, requiring a few megabytes at most. Thus the algorithms typically also assume that the entire data can fit in main memory, and do not address computational issues regarding scalability and “out-of-core” operations.

The tremendous explosion in the amount of data gathering and warehousing in the past few years has generated very large and complex databases. Any effort in mining information from such databases has to address the fact that

- (i) data may be kept in several files as in interlinked relational databases, and information needed for decision making may be spread over more than one file. For example, the concept of “collective data mining” [Kargupta and Park, 2000] explicitly addresses “vertical partitioning” situations where the features or variables relevant to a classification decision are spread over multiple files, each accessible to only one classifier.
- (ii) the files may be spread across several disks or even across different geographical locations, and
- (iii) the statistical quality of data may vary widely. For example the percentage of cases involving financial or health-care fraud varies in different regions, and so does the amount of missing information.

One can argue that by transferring all data to a single warehouse and performing a series of merges and joins, one can get a single (albeit very large), flat file. A traditional algorithm can be used after randomizing and subsampling this file. But in real applications this approach may not be feasible because of the computational, bandwidth and storage costs. In certain cases, it may not even be possible for a variety of practical reasons including security, privacy, proprietary nature of data, need for fault tolerant distribution of data and services, real-time processing requirements, statutory constraints imposed by law, etc. [Prodromidis et al., 2000]. Then there are two options. If the owners of the individual databases are willing to provide high level or summary information/decisions such as local classification estimates, and transmit this information to a central location, then a meta-learner can be applied to the component decisions to come up with a final, composite decision. Note that such high level information not only has reduced storage and bandwidth requirements, but also maintains the privacy of individual records [DuMouchel et al., 1999]. Otherwise one has to resort to a distributed computing

framework such as the emerging field of COllective INtelligence (COIN), wherein techniques are developed such that local and independent computations can still increase a desired global utility function [Wolpert and Tumer, 1999].

The first option leads to several issues reminiscent of studies in decision fusion [Dasarathy, 1994] applied largely to multi-sensor fusion and distributed control problems. It is also related to the theory of ensembles [Sharkey, 1996] and integration of multiple learned models [Chan et al., 1996]. But there are substantial new aspects that need to be addressed in a distributed data mining context.

This chapter is rooted in the ensemble framework and shows how order statistics can be used in the design of a “meta-learner” that examines the outputs of multiple distributed classifiers and provides a final decision. Order statistics is one of the key tools of robust statistics, tailored to handling data with outliers. In a distributed data mining scenario in which there is wide variability among the individual classifiers because of the underlying quality of the local data that they examine, a meta-learner should be able to tolerate a few outlier classifier results. The robust properties of order statistics based approaches such as median filtering and m-estimators [Arnold et al., 1992], have been observed in many disciplines. Thus they are an obvious candidate for meta-learning in such environments.

The next section provides a brief review of the meta-learning framework for classification to put the proposed techniques in perspective. Section 3 summarizes the relationship between classifier errors and decision boundaries and provides the necessary background for mathematically analyzing order statistic combiners [Tumer and Ghosh, 1996a]. Section 4 introduces simple order statistic combiners. Based on these concepts, in Section 5 we propose two powerful combiners, *trim* and *spread*, and derive the amount of error reduction associated with each. In Section 6 we present the performance of order statistic combiners on several datasets. Section 7 discusses the implications of using linear combinations of order statistics as a strategy for pooling the outputs of individual classifiers.

2 A Brief History of Multi-Learner Systems

The idea of integrating multiple models for the *same* problem, has been examined for a long time. The main goal is to obtain a better composite global model, with more accurate and reliable estimates or decisions. Some notable early systems include Selfridge’s Pandemonium [Selfridge, 1958] where a head-deamon would select the daemon that “shouted the loudest”, and Nilsson’s committee machines. A strong motivation for such systems was voiced by Kanal in his classic 1974 paper [Kanal, 1974]:

“It is now recognized that the key to pattern recognition problems does not lie wholly in learning machines, statistical approaches, spatial, filtering,..., or in any other particular solution which has been vigorously advocated by one or another group during the last one and a half decades as the solution to the pattern recognition problem. No single model exists for all pattern recognition problems and no single technique is applicable to all problems. Rather what we have is a bag of tools and a bag of problems.”

In the late seventies, much work was done on combining linguistic and statistical models, and on combining heuristic search with statistical pattern recognition. Subsequently, similar sentiments on the importance of multiple approaches were also voiced in the AI community [Minsky, 1991]:

“To solve really hard problems, we’ll have to use several different representations....It is time to stop arguing over which type of pattern-classification technique is bestInstead we should work at a higher level of organization and discover how to build managerial systems to exploit the different virtues and evade the different limitations of each of these ways of comparing things.”

Integration of multiple data sources and/or learned models can now be found in several disciplines, for example, the combining of estimators in econometrics [Granger, 1989], evidences in rule-based systems [Barnett, 1981] and multi-sensor data fusion [Dasarathy, 1994]. Of course, one can find numerous examples in the human central nervous system [Shepherd, 1979], as well as in some large engineering systems such as those that demand fault tolerance or employing control mechanisms that may need to function in different operating regimes [Narendra et al., 1995]. In particular, multiple models for nonlinear control has a long tradition (see <http://www.itk.ntnu.no/ansatte/Johansen.Tor.Arne/mmamc/address.html> for a detailed list of researchers in this area). Hybridization in a broader sense is seen in efforts to combine two or more of neural network, Bayesian, GA, fuzzy logic and knowledge-based systems [Aggarwal et al., 1996, Taha and Ghosh, 1997]. The goal is again to incorporate diverse sources and forms of information and to exploit the somewhat complementary nature of different methodologies.

Figure 1 shows a generic diagram of an *ensemble*, the simplest and most well understood type of multi-learner systems. Each component learner is a regressor or classifier, trying to solve the same task. While data ultimately originates from an underlying universal set X , each learner may receive somewhat different subsets of the data for “training” or parameter estimation (as in bagging [Breiman, 1994] and boosting [Drucker et al., 1994]), and may be using different feature extractors (f_s) on the same raw data. For example, in our earlier work on sonar classification [Ghosh et al., 1992a], Fourier, wavelet and autoregressive coefficients extracted from the same

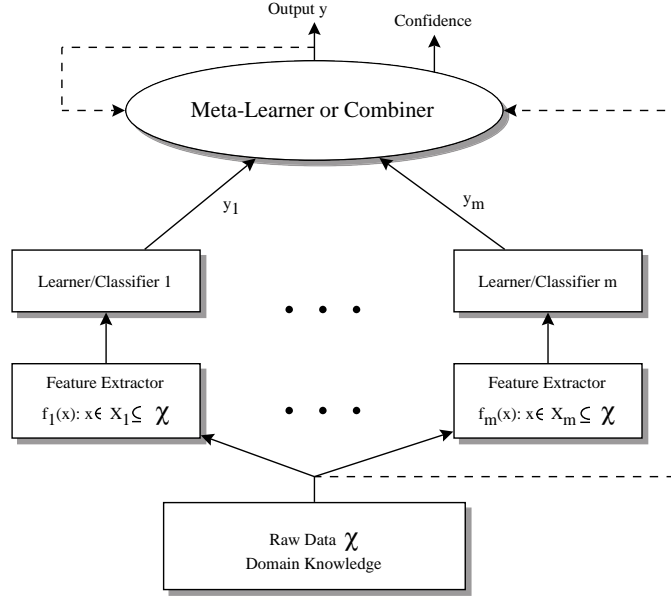


Figure 1: Generic architecture of a multi-learner system for regression or classification

preprocessed time series were used respectively for three different classifier types.

Along with selection of training samples and feature extractors, one needs to decide what types of learners to use and how many, and finally, how to design the meta-learner. There are also larger issues of how to train the components given that they are part of a bigger system, and to estimate the overall gains achievable.

The simplest meta-learner is the *combiner*, where the output y is determined solely from the outputs of the individual learners. In the past few years, a host of experimental results from both neural network and machine learning communities show that combining the outputs of multiple regressors or classifiers via (weighted) averaging, majority vote, product rule, entropy, etc., provides statistically significant improvement in performance along with tighter confidence intervals. Moreover, theoretical analysis has been developed for both regression [Perrone, 1993, Hashem and Schmeiser, 1993] and classification [Tumer and Ghosh, 1996a, Tumer and Ghosh, 1999], to estimate the gains achievable.

Beyond simple combiners are meta-learning methods such as arbitration [Chan and Stolfo, 1997] and stacking [Wolpert, 1992]. Also in this category are divide-and-conquer approaches, where relatively simple learners specialize in different parts of the input-output space, and the total model is a (possibly soft) union of such simpler models. Techniques of modular learning include “mixtures-of-experts”, local linear regression, CART/MARS, adaptive subspace models, etc. [Jordan and Jacobs, 1994, Ramamurti and Ghosh, 1999, Holmstrom et al., 1997]. What

distinguishes all these models from simple combiners is that the meta-learner’s actions now also depends on the current input and/or the target values, i.e. the dotted lines in Fig. 1 are also used.

The simple combining methods are best suited for problems where the individual classifiers perform the same task, and have comparable success. However, such combiners are more susceptible to outliers and to unevenly performing classifiers. On the other hand, the more general meta-learners are conceptually more powerful but are vulnerable to all the problems associated with the added learning (e.g., overparameterizing, lengthy training time). They may also require access to the entire input database which may be impractical in a distributed data mining environment.

The performance of classifier ensembles has been spectacular. Brieman calls the combination of decision tree classifiers with boosting the most significant development in classifier design in this decade. Indeed, ensemble methods are even becoming available in commercial data mining tools such as SAS Enterprise Miner Version 3.0. To see intuitively why ensembles have been so effective, note that different types of classifiers have different “inductive biases” [Geman et al., 1992, Mitchell, 1997], and thus, in general, they do not generalize in identical ways even when they are trained on the same data set, and have comparable performance on a test set. Traditionally, the classifier perceived as the “best”, as indicated by a suitable scoring function such as a cross-validation based *estimation* of the true generalization error, is selected. This means that the other classifiers that had been designed during the model exploration/development phase, get discarded. But this results in a potential loss of useful information and effort. An ensemble can effectively make use of such complementary information to reduce model variance [Perrone, 1993, Tumer and Ghosh, 1999] and in certain situations it also reduces bias as shown by the theory of large margin classifiers [Bartlett and Shawe-Taylor, 1998]. It works best when each learner is well trained, but different learners generalize in different ways, i.e., there is diversity in the ensemble [Krogh and Vedelsby, 1995]. Diversity may be induced through different presentations of the input data, as in bagging, variations in learner design, or by adding a penalty to the outputs to encourage diversity.

For large data sets, there is also a computational reason for using meta-learners, namely, it can make certain inductive learners fairly scalable to large data sets [Chan and Stolfo, 1997, Provost and Kolluri, 1999, Provost, 2000]. In [Chan and Stolfo, 1997] each component classifier uses a *randomly chosen* subset of the entire data set, and the decision tree classifiers can operate in parallel. This proves quite effective as overall computation can be greatly reduced with little loss in performance. A beneficial side-effect of using smaller datasets for decision tree classifiers

is that the classifiers obtained are smaller in size [Oates and Jensen, 1998]. These two as well as other scalability techniques for decision tree based classifiers are covered in a nice recent survey [Provost and Kolluri, 1999].

A fundamental assumption in all the multi-classifier approaches mentioned above is that the designer has access to the entire data set, which can be used in its entirety, resampled in a random (bagging) or weighted (boosting) way [Breiman, 1999], or randomly partitioned and distributed. Thus, except for boosting situations, each classifier sees training data of comparable quality. If the individual classifiers are then appropriately chosen and trained properly, their performances will be (relatively) comparable in any region of the problem space. So gains from combining are derived from the diversity among classifiers rather than by compensating for weak members of the pool.

This assumption is clearly invalid for distributed data mining using heterogeneous sites [Kargupta and Park, 2000]. Such real-life conditions often result in a pool of classifiers that may have significant variations in their overall performance. Moreover, they may lead to conditions where individual classifiers have similar *average* performance, but substantially different performance over different parts of the input space. In such cases, combining is still desirable, but neither simple combiners nor meta-learners are particularly well-suited for the type of problems that arise. For example, the simplicity of averaging the classifier outputs is appealing, but the prospect of one poor classifier corrupting the combiner makes this a risky choice. Weighted averaging of classifier outputs appears to provide some flexibility [Hashem and Schmeiser, 1993, Merz and Pazzani, 1997]. Unfortunately, the weights are still assigned on a per classifier basis rather than a per sample or per class basis. If a classifier is accurate only in certain areas of the input space, this scheme fails to take advantage of the variable accuracy of the classifier in question. Using a meta-learner that provides different weights for different patterns can potentially solve this problem, but at a considerable cost. Also, as explained earlier, the off-line training of a meta-learner using substantial amount of data outputted by geographically distributed classifiers, may not be feasible or even allowable. In addition to providing robustness demanded in such situations, the order statistic combiners presented in this work also aim at bridging the gap between simplicity and generality by allowing the flexible selection of classifiers without the associated cost of training meta-classifiers.

3 Error Characterization in a Single Classifier

In this section we summarize the approach and results of [Tumer and Ghosh, 1996a, Tumer and Ghosh, 1999]¹, that quantify the effect of inaccuracies in estimating *a posteriori* class probabilities on the classification error for a *single* classifier. This background is needed to characterize and understand the impact of order statistics combinners, as described in Sections 3 and 4.

It is well known that, given *one-of-L* desired outputs and sufficient training samples reflecting the class priors, the outputs of certain classifiers trained to minimize a mean square or cross-entropy error criteria, approximate the *a posteriori* probability densities of the corresponding classes [Richard and Lippmann, 1991, Ruck et al., 1990]. Based on this result, one can model the *i*th output of the *m*th such classifier as:

$$f_i^m(x) = p_i(x) + \epsilon_i^m(x), \quad (1)$$

where $p_i(x)$ is the true posterior for *i*th class on input x , and $\epsilon_i^m(x)$ is the error of the *m*th classifier in estimating that posterior.

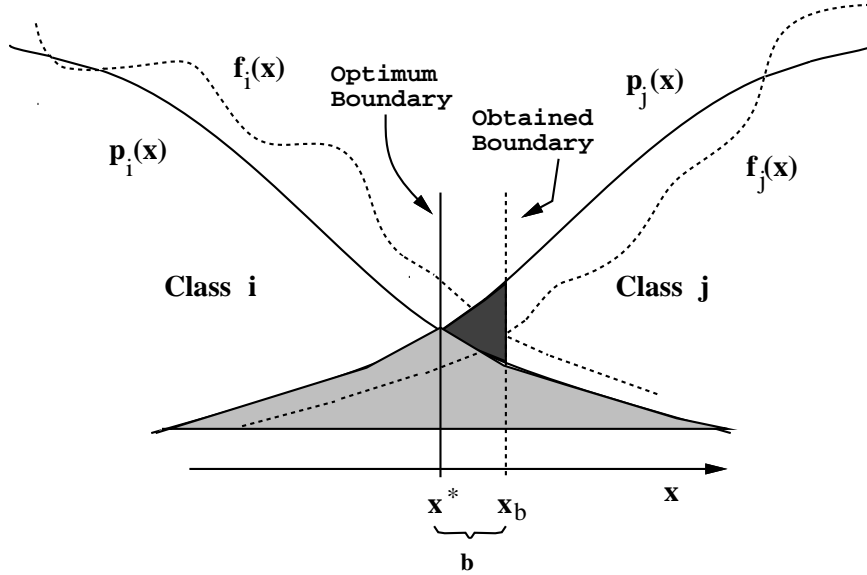


Figure 2: Error regions associated with approximating the *a posteriori* probabilities [Tumer and Ghosh, 1996a].

Now, let us decompose the error into two parts: $\epsilon_i^m(x) = \beta_i^m + \eta_i^m(x)$. The first component does not vary with the input, and provides an offset, or systematic error for each class. The

¹This and other related papers can be downloaded from URL <http://www.lans.ece.utexas.edu>.

second component gives the variability from that systematic error, for each x in each class, and has zero mean and variance $\sigma_{\eta_i^m(x)}^2$. These two components of the error are similar to the bias and variance decomposition for a quadratic loss function given in [Geman et al., 1992], although they are at the individual input level. We will therefore refer to classifiers as “biased” and “unbiased” implying $\beta_k^m \neq 0$ for some k, m , and $\beta_k^m = 0$, $\forall k, m$, respectively. Let b^m denote the offset between the ideal class boundary, x^* (based on $p_i(x) = p_j(x)$) and the realized boundary, x_b^m (based on $f_i^m(x) = f_j^m(x)$), as shown in Figure 2 [Tumer and Ghosh, 1996a]. This boundary offset ($b^m = x_b^m - x^*$) has mean and variance given respectively by:

$$\beta^m = \frac{\beta_i^m - \beta_j^m}{s}, \quad (2)$$

and

$$\sigma_{b^m}^2 = \frac{\sigma_{\eta_i^m(x)}^2 + \sigma_{\eta_j^m(x)}^2}{s^2}, \quad (3)$$

where $s = p'_j(x^*) - p'_i(x^*)$ as introduced in [Tumer and Ghosh, 1996a].

Let us further denote the probability density function of this boundary offset by $f_b(x)$. The expected model error associated with the selection of a particular classifier m , can then be expressed as:

$$E_{model}^m = \int_{-\infty}^{\infty} A(b) f_b(b) db, \quad (4)$$

where $A(b) = \int_{x^*}^{x^*+b} (p_j(x) - p_i(x)) dx$ is the error due to the selection of a particular decision boundary. In general, it is not possible to obtain the density function for the boundary offset without making assumptions on the distributions of the errors. However, a first order approximation, derived in [Tumer and Ghosh, 1996a], leads to:

$$E_{model}^m = \int_{-\infty}^{\infty} \frac{1}{2} b^2 s f_b(b) db. \quad (5)$$

Let us define the first and second moments of the boundary offset as follows:

$$\mathcal{M}_1 = \int_{-\infty}^{\infty} x f_b(x) dx \quad \text{and} \quad \mathcal{M}_2 = \int_{-\infty}^{\infty} x^2 f_b(x) dx.$$

If the individual classifiers are unbiased, the offset b^m of a single classifier has $\mathcal{M}_1 = 0$ and $\mathcal{M}_2 = \sigma_{b^m}^2$, leading to:

$$E_{model}^m = \frac{s \mathcal{M}_2}{2} = \frac{s \sigma_{b^m}^2}{2}. \quad (6)$$

Now, if the classifiers are biased, the variance of b is left unchanged (given by Equation 3), but the mean becomes $\beta = \frac{\beta_i - \beta_j}{s}$. In other words, we have $\mathcal{M}_1 = \beta^m$ and $\sigma_{b^m}^2 = \mathcal{M}_2 - \mathcal{M}_1^2$, leading to the following model error:

$$E_{model}^m(\beta) = \frac{s \mathcal{M}_2}{2} = \frac{s}{2} (\sigma_{b^m}^2 + (\beta^m)^2). \quad (7)$$

To emphasize the distinction between biased and unbiased classifiers, the model error will be given as a function of β for biased classifiers. A more detailed derivation of class boundaries and error regions is presented in [Tumer and Ghosh, 1996a]. For analyzing the error regions after combining and comparing them to the single classifier case, one needs to determine how the first and second moments of the boundary distributions are affected by combining. The following sections focus on obtaining those values for order statistics based combiners.

4 Combining Multiple Classifiers through Order Statistics

4.1 Basic Concepts

In this section, we briefly discuss some basic concepts and properties of order statistics. Let X be a random variable with probability density function $f_X(\cdot)$, and cumulative distribution function $F_X(\cdot)$. Let (X_1, X_2, \dots, X_N) be a random sample drawn from this distribution. Now, let us arrange them in non-decreasing order, providing:

$$X_{1:N} \leq X_{2:N} \leq \dots \leq X_{N:N}.$$

The i th order statistic denoted by $X_{i:N}$, is the i th value in this progression. The cumulative distribution function for the smallest and largest order statistic can be obtained by noting that:

$$F_{X_{N:N}}(x) = P(X_{N:N} \leq x) = \prod_{i=1}^N P(X_{i:N} \leq x) = [F_X(x)]^N$$

and:

$$\begin{aligned} F_{X_{1:N}}(x) &= P(X_{1:N} \leq x) = 1 - P(X_{1:N} \geq x) = 1 - \prod_{i=1}^N P(X_{i:N} \geq x) \\ &= 1 - (1 - \prod_{i=1}^N P(X_{i:N} \leq x)) = 1 - [1 - F_X(x)]^N \end{aligned}$$

The corresponding probability density functions can be obtained from these equations. In general, for the i th order statistic, the cumulative distribution function gives the probability that exactly i of the chosen X 's are less than or equal to x . The probability density function of $X_{i:N}$ is then given by [David, 1970]:

$$f_{X_{i:N}}(x) = \frac{N!}{(i-1)!(N-i)!} [F_X(x)]^{i-1} [1 - F_X(x)]^{N-i} f_X(x). \quad (8)$$

This general form however, cannot always be computed in closed form. Therefore, obtaining the expected value of a function of x using Equation 8 is not always possible. However, the first two moments of the density function are widely available for a variety of distributions [Arnold et al., 1992]. These moments can be used to compute the expected values of certain specific functions, e.g., polynomials of order less than two.

4.2 Combining Unbiased Classifiers through Order Statistics

Now, let us turn our attention to order statistics (OS) combiners. For a given input x , let the network outputs of each of the N classifiers for each class i be ordered in the following manner:

$$f_i^{1:N}(x) \leq f_i^{2:N}(x) \leq \dots \leq f_i^{N:N}(x).$$

Then one constructs the k th order statistic combiner, by selecting the k th ranked output for each class ($f_i^{k:N}(x)$), as representing its posterior.

In particular, *max*, *med* and *min* combiners are defined as follows:

$$f_i^{max}(x) = f_i^{N:N}(x), \quad (9)$$

$$f_i^{med}(x) = \begin{cases} \frac{f_i^{\frac{N}{2}:N}(x) + f_i^{\frac{N}{2}+1:N}(x)}{2} & \text{if } N \text{ is even} \\ f_i^{\frac{N+1}{2}:N}(x) & \text{if } N \text{ is odd,} \end{cases} \quad (10)$$

$$f_i^{min}(x) = f_i^{1:N}(x). \quad (11)$$

These three combiners are relevant because they represent important qualitative interpretations of the output space. Selecting the maximum combiner is equivalent to selecting the class with the highest posterior. Indeed, since the network outputs approximate the class *a posteriori* distributions, selecting the maximum reduces to selecting the classifier that is the most “certain” of its decision. The drawback of this method however is that it can be compromised by a single classifier that repeatedly provides high values. The selection of the minimum combiner follows a similar logic, but focuses on classes that are unlikely to be correct, rather than on the correct class. Thus, this combiner eliminates less likely classes by basing the decision on the lowest value for a given class. This combiner suffers from the same ills as the *max* combiner. However, it is less dependent on a single error, since it performs a min-max operation, rather than a max-max². The median classifier on the other hand considers the most “typical” representation of each class. For highly noisy data, this combiner is more desirable than either the *min* or *max* combiners since the decision is not compromised as much by a single large error.

The analysis that follows does not depend on the particular order statistic chosen. Therefore, we will denote all OS combiners by $f_k^{os}(x)$ and derive the model error, E_{model}^{os} . The network output provided by $f_k^{os}(x)$ is given by:

$$f_k^{os}(x) = p_k(x) + \epsilon_k^{os}(x), \quad (12)$$

Let us first investigate the zero-bias case ($\beta_k = 0$, $\forall k$), where we get $\epsilon_k^{os}(x) = \eta_k^{os}(x)$.

²Recall that the pattern is ultimately assigned to the class with the highest combined output.

Proceeding as in Section 3, the boundary b^{os} is shown to be:

$$b^{os} = \frac{\eta_i^{os}(x_b) - \eta_j^{os}(x_b)}{s}. \quad (13)$$

For *i.i.d.* η_k 's, the first two moments will be identical for each class. Moreover, taking the order statistic will shift the mean of both η_i^{os} and η_j^{os} by the same amount, leaving the mean of the difference unaffected. Therefore, b^{os} will have zero mean, and variance:

$$\sigma_{b^{os}}^2 = \frac{2 \sigma_{\eta_k^{os}}^2}{s^2} = \frac{2 \alpha \sigma_{\eta_k^m}^2}{s^2} = \alpha \sigma_{b^m}^2, \quad (14)$$

where α is a reduction factor that depends on the order statistic and on the distribution of b . For most distributions, α can be found in tabulated form [Arnold et al., 1992]. For example, Table 1 provides α values for all order statistic combiners, up to 10 classifiers, for a Gaussian distribution [Arnold et al., 1992, Sarhan and Greenberg, 1956]. (Because this distribution is symmetric, the α values of l and k where $l + k = N + 1$ are identical, and listed in parenthesis).

Returning to the error calculation, we have: $\mathcal{M}_1^{os} = 0$, and $\mathcal{M}_2^{os} = \sigma_{b^{os}}^2$, providing:

$$E_{model}^{os} = \frac{s \mathcal{M}_2^{os}}{2} = \frac{s \sigma_{b^{os}}^2}{2} = \frac{s \alpha \sigma_{b^m}^2}{2} = \alpha E_{model}^m. \quad (15)$$

Table 1: Reduction factors α for the Gaussian Distribution, based on [Sarhan and Greenberg, 1956].

N	k	α	N	k	α	N	k	α
1	1	1.00	6	2 (5)	.280	1	(9)	.357
2	1 (2)	.682	3	(4)	.246	2	(8)	.226
3	1 (3)	.560	1	(7)	.392	9	3 (7)	.186
	2	.449	7	2 (6)	.257	4	(6)	.171
4	1 (4)	.492	3	(5)	.220	5		.166
	2 (3)	.360	4		.210	1	(10)	.344
	1 (5)	.448	1	(8)	.373	2	(9)	.215
5	2 (4)	.312	8	2 (7)	.239	10	3 (8)	.175
	3	.287	3	(6)	.201	4	(7)	.158
6	1 (6)	.416	4	(5)	.187	5	(6)	.151

Equation 15 shows that the reduction in the error due to using the OS combiner instead of the m th classifier is directly related to the reduction in the variance of the boundary offset b . Since the means and variances of order statistics for a variety of distributions are widely available in tabular form, the reductions can be readily quantified.

4.3 Combining Biased Classifiers through Order Statistics

In this section, we analyze the error regions for biased classifiers. Let us return our attention to b^{os} . First, note that the error terms can no longer be studied separately, since in general $(a + b)^{os} \neq a^{os} + b^{os}$. We will therefore need to specify the mean and variance of the result of each operation³. Equation 13 becomes:

$$b^{os} = \frac{(\beta_i + \eta_i(x_b))^{os} - (\beta_j + \eta_j(x_b))^{os}}{s}. \quad (16)$$

Let $\bar{\beta}_k = \frac{1}{N} \sum_{m=1}^N \beta_k^m$ be the mean of classifier biases. Since η_k^m 's have zero-mean, $\beta_k + \eta_k(x_b)$ has first moment $\bar{\beta}_k$ and variance $\sigma_{\eta_k^m}^2 + \sigma_{\beta_k^m}^2$, with $\sigma_{\beta_k^m}^2 = E[(\beta_k^m)^2] - \bar{\beta}_k^2$, where $[\cdot]$ denotes the expected value operator.

Taking a specific order statistic of this expression will modify both moments. The first moment is given by $\bar{\beta}_k + \mu^{os}$, where μ^{os} is a shift which depends on the order statistic chosen, but not on the class. Then, the first moment of b^{os} is given by:

$$\frac{(\bar{\beta}_i + \mu^{os}) - (\bar{\beta}_j + \mu^{os})}{s} = \frac{\bar{\beta}_i - \bar{\beta}_j}{s} = \bar{\beta}. \quad (17)$$

Note that the bias term represents an ‘‘average bias’’ since the contributions due to the order statistic are removed. Therefore, reductions in bias cannot be obtained from a table similar to Table 1.

Now, let us turn our attention to the variance. Since $\beta_k^m + \eta_k^m(x_b)$ has variance $\sigma_{\eta_k^m}^2 + \sigma_{\beta_k^m}^2$, it follows that $(\beta_k + \eta_k(x_b))^{os}$ has variance $\sigma_{\eta_k^{os}}^2 = \alpha(\sigma_{\eta_k^m}^2 + \sigma_{\beta_k^m}^2)$, where α is the factor discussed in Section 4.2. Therefore, the variance of b^{os} is given by:

$$\begin{aligned} \sigma_{b^{os}}^2 &= \frac{\sigma_{\eta_i^{os}}^2 + \sigma_{\eta_j^{os}}^2}{s^2} = \frac{2 \alpha \sigma_{\eta_i^m}^2}{s^2} + \frac{\alpha(\sigma_{\beta_i^m}^2 + \sigma_{\beta_j^m}^2)}{s^2} \\ &= \alpha(\sigma_{b^m}^2 + \sigma_{\beta^m}^2), \end{aligned} \quad (18)$$

where $\sigma_{\beta^m}^2 = \frac{\sigma_{\beta_i^m}^2 + \sigma_{\beta_j^m}^2}{s^2}$ is the variance introduced by the systematic errors of different classifiers.

We have now obtained the first and second moments of b^{os} , and can compute the model error. Namely, we have $\mathcal{M}_1^{os} = \bar{\beta}$ and $\sigma_{b^{os}}^2 = \mathcal{M}_2^{os} - (\mathcal{M}_1^{os})^2$, leading to:

$$E_{model}^{os}(\beta) = \frac{s}{2} \mathcal{M}_2^{os} = \frac{s}{2} (\sigma_{b^{os}}^2 + \bar{\beta}^2) \quad (19)$$

$$= \frac{s}{2} (\alpha(\sigma_{b^m}^2 + \sigma_{\beta^m}^2) + \bar{\beta}^2). \quad (20)$$

The reduction in the error is more difficult to assess in this case. By writing the error as:

$$E_{model}^{os}(\beta) = \alpha \frac{s}{2} (\sigma_b^2 + (\beta^m)^2) + \frac{s}{2} (\alpha \sigma_\beta^2 + \bar{\beta}^2 - \alpha (\beta^m)^2),$$

³Since the exact distribution parameters of b^{os} are not known, we use the sample mean and the sample variance.

we get:

$$E_{model}^{os}(\beta) = \alpha E_{model}^m(\beta) + \frac{s}{2} (\alpha \sigma_{\beta}^2 + \bar{\beta}^2 - \alpha(\beta^m)^2). \quad (21)$$

Analyzing the error reduction in the general case requires knowledge about the bias introduced by each classifier. Unlike regression problems where the bias and variance contributions to the error are additive and well-understood, in classification problems their interaction is more complex [Friedman, 1997]. Indeed it has been observed that ensemble methods do more than simply reduce the variance [Schapire et al., 1997].

Based on these observations and Equation 21, let us analyze extreme cases. For example, if each classifier has the same bias, σ_{β}^2 is reduced to zero and $\bar{\beta} = \beta^m$. In this case the error reduction can be expressed as:

$$E_{model}^{os}(\beta) = \frac{s}{2} (\alpha \sigma_b^2 + (\beta^m)^2) = \alpha E_{model}^m(\beta) + \frac{s(1-\alpha)}{2} (\beta^m)^2,$$

where α balances the two contributions to the error. A small value for α will reduce the first component of the error (mainly variance), while leaving the second term untouched. The net effect will be very similar to results obtained for regression problems. In this case, it is important to reduce classifier bias before combining (e.g., by using an overparametrized model).

If on the other hand, the biases produce a zero mean variable, we obtain $\bar{\beta} = 0$. In this case, the model error becomes:

$$E_{model}^{os}(\beta) = \alpha E_{model}^m(\beta) + \frac{s\alpha}{2} (\sigma_{\beta^m}^2 - (\beta^m)^2)$$

and the error reduction will be significant if the second term is small or negative. In fact, if the variation among the biases is small relative to their magnitude, the error will be reduced more than in the unbiased cases. If however, the variation is large compared to the magnitude, the error reduction will be minimal. Furthermore, if α is large and the biases are small and highly varied, it is possible for this combiner to do worse than the individual classifiers, which is a danger not present for regression problems. This observation very closely parallels results reported in [Friedman, 1997].

5 Linear Combining of Ordered Classifier Outputs

In the previous section, we derived error reductions when the class posteriors are directly estimated through the ordered classifier outputs. Since simple averaging has also been shown to provide benefits, in this section, we investigate the combinations of averaging and order statistics for pooling classifier outputs.

5.1 Spread Combiner

The first linear combination of ordered classifier outputs we study focuses on extrema. As discussed in Section 4.2, the maximum and minimum of a set of classifier outputs carry specific meanings. Indeed, the maximum can be viewed as the class for which there is the most evidence. Similarly, the minimum deletes classes with little evidence. In order to avoid a single classifier from having too large of an impact on the eventual output, these two values can be averaged to yield the *spread* combiner. This combiner strikes a balance between the positive and negative evidence, leading to a more robust combiner than either of them.

5.1.1 Spread Combiner for Unbiased Classifiers:

For a classifier without bias, the spread combiner is formally defined as:

$$f_i^{spr}(x) = \frac{1}{2} (f_i^{1:N}(x) + f_i^{N:N}(x)) = p(c_i|x) + \eta_i^{spr}(x), \quad (22)$$

where:

$$\eta_i^{spr}(x) = \frac{1}{2} (\eta_i^{1:N}(x) + \eta_i^{N:N}(x)).$$

The variance of $\eta_i^{spr}(x)$ is given by:

$$\sigma_{\eta_i^{spr}}^2 = \frac{1}{4} \sigma_{\eta_i^{1:N}}^2 + \frac{1}{4} \sigma_{\eta_i^{N:N}}^2 + \frac{1}{2} \text{cov}(\eta_i^{1:N}(x), \eta_i^{N:N}(x)). \quad (23)$$

where $\text{cov}(\cdot, \cdot)$ represents the covariance between two variables (even when the η_i 's are independent, ordering introduces correlations). Note that because of the ordering, the variances in the first two terms of Equation 23 can be expressed in terms of the individual classifier variances. Furthermore, the covariance between two order statistics can also be determined in tabulated form for given distributions. Table 2 provides these values for a Gaussian distribution based on [Sarhan and Greenberg, 1956]. This expression can be further simplified for symmetric distributions where $\sigma_{\eta^{1:N}}^2 = \sigma_{\eta^{N:N}}^2$ (e.g., Gaussian noise model) and leads to:

$$\sigma_{\eta_i^{spr}}^2 = \frac{1}{2} (\alpha_{1:N} + B_{1,N:N}) \sigma_{\eta_i(x)}^2, \quad (24)$$

where $\alpha_{m:N}$ is the variance of the m th ordered sample and $B_{m,l:N}$ is the covariance between the m th and l th ordered samples, given that the initial samples had unit variance [Sarhan and Greenberg, 1956]. Because this is a symmetric distribution, the β values are also symmetric (e.g., $\beta_{1,2:5} = \beta_{4,5:5}$).

Then, using Equation 3, the variance of the boundary offset b^{spr} can be calculated:

$$\begin{aligned} \sigma_{b^{spr}}^2 &= \frac{\sigma_{\eta_i^{spr}}^2 + \sigma_{\eta_j^{spr}}^2}{s^2} \\ &= \frac{1}{2} (\alpha_{1:N} + B_{1,N:N}) \sigma_b^2. \end{aligned} \quad (25)$$

Table 2: Some Reduction Factors B for the Gaussian Distribution, based on [Sarhan and Greenberg, 1956].

N	k, l	B	N	k, l	B	N	k, l	B	N	k, l	B	
2	1,2	.318	6	2,3	.189	8	1,4	.095	9	1,6	.059	
3	1,2	.276		2,4	.140		1,5	.075		1,7	.049	
	1,3	.165		2,5	.106		1,6	.060		1,8	.040	
	1,2	.246		3,4	.183		1,7	.048		1,9	.031	
4	1,3	.158	7	1,2	.196		1,8	.037		2,3	.154	
	1,4	.105		1,3	.132		2,3	.163		2,4	.117	
	2,3	.236		1,4	.099		2,4	.123		2,5	.093	
	1,2	.224		1,5	.077		2,5	.098		2,6	.077	
	1,3	.148		1,6	.060		2,6	.079			2,7	.063
5	1,4	.106	7	1,7	.045		2,7	.063			2,8	.052
	1,5	.074		2,3	.175		3,4	.152			3,4	.142
	2,3	.208		2,4	.131		3,5	.121			3,5	.114
	2,4	.150		2,5	.102		3,6	.098			3,6	.093
	1,2	.209		2,6	.080		4,5	.149			3,7	.077
	1,3	.139		3,4	.166		1,2	.178			4,5	.137
6	1,4	.102		3,5	.130		1,3	.121			4,6	.113
	1,5	.077	8	1,2	.186		9	1,4	.091			
	1,6	.056		1,3	.126		1,5	.073				

Finally, through Equation 6, we can obtain the reduction in the model error due to the spread combiner:

$$\frac{E_{model}^{spr}}{E_{model}} = \frac{\alpha_{1:N} + B_{1,N:N}}{2} . \quad (26)$$

Based on Equation 26 and Tables 1 and 2, Table 3 displays the error reductions provided by the spread combiner for a Gaussian noise model (for comparison purposes, the error reduction for the *min* and *max* combiners is also provided. Note that for the Gaussian distribution, the error reduction of *min* is equal to that of *max*).

Table 3: Error Reduction Factors for the Spread, *min* and *max* Combiners with Gaussian Noise Model.

N	<i>spread</i>	<i>min</i> or <i>max</i>
2	.500	.682
3	.362	.560
4	.299	.492
5	.261	.448
6	.236	.416
7	.219	.392
8	.205	.373
9	.194	.357
10	.186	.344

5.1.2 Spread Combiner for Biased Classifiers:

Now, if the classifier biases are non-zero, the spread combiner's output is given by:

$$f_i^{spr}(x) = \frac{1}{2} (f_i^{1:N}(x) + f_i^{N:N}(x)) = p(c_i|x) + (\eta_i(x) + \beta_i)^{spr}. \quad (27)$$

In that case, the boundary offset is given by:

$$b^{spr} = \frac{(\beta_i + \eta_i(x_b))^{spr} - (\beta_j + \eta_j(x_b))^{spr}}{s}, \quad (28)$$

which after expanding each term and regrouping can be expressed as:

$$b^{spr} = \frac{(\beta_i + \eta_i(x_b))^{1:N} - (\beta_j + \eta_j(x_b))^{1:N}}{2s} + \frac{(\beta_i + \eta_i(x_b))^{N:N} - (\beta_j + \eta_j(x_b))^{N:N}}{2s}. \quad (29)$$

The first moment of b^{spr} can be obtained by analyzing each term of Equation 29. In fact, the offset introduced by the first and n th order statistic for classes i and j will cancel each other out, leaving only the average bias between the min and max components of the error (as in Equation 17), given by $\beta^{spr} = \frac{\beta_i^{1:N} - \beta_j^{1:N} + \beta_i^{N:N} - \beta_j^{N:N}}{s}$.

The variance of b^{spr} needs to be derived from Equation 29. Proceeding as in Equation 18, the variance of the spread combiner can be expressed as:

$$\sigma_{b^{spr}}^2 = \left(\frac{1}{4}\alpha_{1:N} + \frac{1}{4}\alpha_{N:N} + \frac{1}{2}B_{1,N:N} \right) (\sigma_{b_m}^2 + \sigma_{\beta_m}^2). \quad (30)$$

For a symmetric distribution (where $\alpha_{1:N} = \alpha_{N:N}$), we obtain the following error:

$$\begin{aligned}
E_{model}^{spr}(\beta) &= \frac{s}{2} \mathcal{M}_2 = \frac{s}{2} (\sigma_{b^{spr}}^2 + \mathcal{M}_1^2) \\
&= \frac{s}{2} \left(\frac{1}{2} \alpha_{1:N} + \frac{1}{2} B_{1,N:N} \right) (\sigma_{b^m}^2 + \sigma_{\beta^m}^2) + (\beta^{spr})^2 \\
&= \frac{1}{2} (\alpha_{1:N} + B_{1,N:N}) E_{model}(\beta) + \\
&\quad \frac{s}{4} (\alpha_{1:N} + B_{1,N:N}) (\sigma_{\beta^m}^2 - (\beta^m)^2) + \frac{s}{2} (\beta^{spr})^2, \tag{31}
\end{aligned}$$

which is very similar to Equation 21, where the value of α for a single order statistic is now replaced by $\frac{\alpha_{1:N} + B_{1,N:N}}{2}$, since the mean of the first and n th order statistic is used in the posterior estimate.

5.2 Trimmed Means

Instead of actively using the extreme values as was the case with the spread combiner, one can base the posterior estimate around the median values. However, instead of selecting one classifier output as was done for f^{med} , one can use multiple classifiers whose outputs are “typical.” In this scheme, only a certain fraction of all available classifiers are used *for a given* pattern. The main advantage of this method over weighted averaging is that the set of classifiers which contribute to the combiner vary from pattern to pattern. Furthermore, they do not need to be determined externally, but are a function of the current pattern and the classifier responses to that pattern.

5.2.1 Trimmed Mean Combiner for Unbiased Classifiers:

Let us formally define the trimmed mean combiner ($\beta_k = 0, \forall k$) as follows:

$$f_i^{trim}(x) = \frac{1}{N_2 - N_1 + 1} \sum_{m=N_1}^{N_2} f_i^{m:N}(x) = p(c_i|x) + \eta_i^{trim}(x), \tag{32}$$

where:

$$\eta_i^{trim}(x) = \frac{1}{N_2 - N_1 + 1} \sum_{m=N_1}^{N_2} \eta_i^m(x).$$

The variance of $\eta_i^{trim}(x)$ is given by:

$$\begin{aligned}
\sigma_{\eta_i^{trim}}^2 &= \frac{1}{(N_2 - N_1 + 1)^2} \sum_{l=N_1}^{N_2} \sum_{m=N_1}^{N_2} cov(\eta_i^{m:N}(x), \eta_i^{l:N}(x)) \\
&= \frac{1}{(N_2 - N_1 + 1)^2} \left(\sum_{m=N_1}^{N_2} \sigma_{\eta_i^{m:N}(x)}^2 + \sum_{m=N_1}^{N_2} \sum_{l>m}^{N_2} 2 cov(\eta_i^{m:N}(x), \eta_i^{l:N}(x)) \right). \tag{33}
\end{aligned}$$

Again, using the factors in Tables 1 and 2, Equation 33 can be further simplified. Note that because the Gaussian distribution is symmetric, the covariance between the k th and l th ordered samples is the same as that between the $N + 1 - k$ th and $N + 1 - l$ th ordered samples. Therefore, Equation 33 leads to:

$$\begin{aligned}\sigma_{\eta_i}^{2\text{trim}} &= \frac{1}{(N_2 - N_1 + 1)^2} \sum_{m=N_1}^{N_2} \alpha_{m:N} \sigma_{\eta_i(x)}^2 \\ &+ \frac{2}{(N_2 - N_1 + 1)^2} \sum_{m=N_1}^{N_2} \sum_{l>m} B_{m,l:N} \sigma_{\eta_i(x)}^2 ,\end{aligned}\quad (34)$$

where $\alpha_{m:N}$ is the variance of the m th ordered sample and $B_{m,l:N}$ is the covariance between the m th and l th ordered samples, given that the initial samples had unit variance [Sarhan and Greenberg, 1956]. Using the theory highlighted in Section 3, and Equation 34, we obtain the following model error reduction:

$$\frac{E_{model}^{\text{trim}}}{E_{model}} = \frac{1}{(N_2 - N_1 + 1)^2} \left(\sum_{m=N_1}^{N_2} \alpha_{m:N} + 2 \sum_{m=N_1}^{N_2} \sum_{l>m} B_{m,l:N} \right). \quad (35)$$

Based on Equation 35 and Tables 1 and 2, we have generated a sample *trim* combiner reduction table. Because there are many possibilities for N_1 and N_2 , a table that exhaustively provides all reduction values is not practical. In this sample table we have selected $N_1 = 2$ and $N_2 = N - 1$, that is, averaging after the lowest and highest values have been removed. For comparison purposes the reduction factors of the averaging combiner for N and $N - 2$ classifiers are also provided (for i.i.d. classifiers the reduction factors are $1/N$ as derived in [Tumer and Ghosh, 1996a]; similar results were obtained for regression problems [Perrone and Cooper, 1993]). As these numbers demonstrate, although $N - 2$ classifiers are used in the trim combiner, *selectively* weeding out undesirable classifiers provides reduction factors significantly better than simply averaging $N - 2$ arbitrary classifiers. The *trim* combiner provides reduction factors comparable the the N classifier *ave* combiner without being susceptible to corruption by one particularly faulty classifier.

5.2.2 Trimmed mean Combiner for Biased Classifiers:

Now, if the classifier biases are non-zero, the trimmed mean combiner's output is given by:

$$f_i^{\text{trim}}(x) = \frac{1}{N_2 - N_1 + 1} \sum_{m=N_1}^{N_2} f_i^{m:N}(x) = p(c_i|x) + (\eta_i(x) + \beta_i)^{\text{trim}}. \quad (36)$$

In that case the boundary offset is given by:

$$b^{\text{trim}} = \frac{(\beta_i + \eta_i(x_b))^{\text{trim}} - (\beta_j + \eta_j(x_b))^{\text{trim}}}{s}. \quad (37)$$

Table 4: Error Reduction Factors for Trim and two corresponding *ave* Combiners with Gaussian Noise Model.

N	<i>ave</i> (for N)	<i>trim</i> (for $N_1 = 2$; $N_2 = N - 1$)	<i>ave</i> (for $N - 2$)
3	.333	.449	1.00
4	.250	.298	.500
5	.200	.227	.333
6	.167	.184	.250
7	.143	.155	.200
8	.125	.134	.167
9	.111	.113	.143

The first moment of b^{trim} can be obtained from a manner similar to that of the spread combiner. Indeed, each mean offset introduced by a specific order statistic for class i will be offset by the one introduced for class j . Only the trimmed mean of the biases will remain, giving the first moment of b^{trim} :

$$\beta^{trim} = \frac{1}{N_2 - N_1 + 1} \sum_{m=N_1}^{N_2} \frac{\beta_i^{m:N} - \beta_j^{m:N}}{s}. \quad (38)$$

In deriving the variance of b^{trim} , we follow the same steps as in Sections 4.3 and 5.1.1. The resulting boundary variance is similar to Equation 18, but the since the reduction is due to the linear combination of multiple ordered outputs, α is replaced by \mathcal{A} , where:

$$\mathcal{A} = \frac{1}{(N_2 - N_1 + 1)^2} \left(\sum_{m=N_1}^{N_2} \alpha_{m:N} + 2 \sum_{m=N_1}^{N_2} \sum_{l>m} B_{m,l:N} \right). \quad (39)$$

The model error reduction in this case is given by:

$$\begin{aligned} E_{model}^{trim}(\beta) &= \frac{s}{2} \mathcal{M}_2 = \frac{s}{2} (\sigma_{b^{trim}}^2 + \mathcal{M}_1^2) \\ &= \frac{s}{2} (\mathcal{A} (\sigma_{\beta^m}^2 + \sigma_{\beta^m}^2) + (\beta^{spr})^2) \\ &= \mathcal{A} E_{model}(\beta) + \frac{s}{2} (\mathcal{A} (\sigma_{\beta^m}^2 - (\beta^m)^2) + (\beta^{spr})^2). \end{aligned} \quad (40)$$

Once again we need to look at the interaction between the two parts of the error reduction. The first term provides the error reduction compared to the model error of an individual classifier. The smaller \mathcal{A} is, the more error reduction there will be. In the second term, on the other hand, a small value for \mathcal{A} is only useful if the variability in the individual biases is higher than the biases themselves ($\sigma_{\beta^m}^2 > (\beta^m)^2$).

6 Experimental Results

The order statistics-based combining methods proposed in this article are tailored for situations where one or more of the following apply:

1. Individual classifier performance is uneven and class dependent;
2. It is not possible (insufficient data, high amount of noise) to fine tune the individual classifiers without using computationally expensive methods;
3. All the features may not be available to all the classifiers.

Such situations occur, for example, in electrical logging while drilling for oil, where data from certain well sites almost completely misses out on portions of the problem space, and in imaging from airborne platforms where the classifiers receive inputs from different satellites and/or different types of sensors (e.g., thermal, optical, SAR). In this article we restrict ourselves to public domain datasets and simulate such variability in two ways, namely, by

- (i) segmenting the feature set and allowing individual classifiers to have access to only a limited portion of the feature set.
- (ii) using “early stopping” i.e., prematurely terminating the training of the individual classifiers⁴.

For the experiments reported below, we used a multi-layer perceptron (MLP) with a single hidden layer, whose weights were randomly initialized for each run. All classification results reported in this article are *test set error rates averaged over 20 runs, along with the differences in the mean (standard deviation divided by square root of the number of runs)*. Several types of simple combiners such as averaging, weighted averaging, voting, median, products, weighted products (Bayesian), using Dempster-Schafer theory of evidence, and entropy-based averaging, have been proposed in the literature. However, on a wide variety of data sets, it has been observed that simple averaging usually provides results comparable to any of these techniques (and, surprisingly, often better than most of them) [Ghosh et al., 1992b, Tumer and Ghosh, 1996b]. For this reason, in this study, we use the average combiner as a *representative* of simple combiners, for comparison purposes.

6.1 Variability through Segmentation

The first group of experiments focus on classifiers that because of circumstances (e.g., geography) have access to only a part of the full feature set. This situation fits the collective data mining framework wherein a globally distributed dataset is vertically partitioned

⁴In all the experiments reported here, “high variability” means that classifiers in an ensemble were trained half as long as they would have been, had they been stand-alone classifiers.

[Kargupta and Park, 2000]. Unfortunately, we are not aware of any public domain datasets for collective data mining, so we instead selected three data sets from the Proben1/UCI benchmarks [Prechelt, 1994]. Briefly these data sets, and the corresponding size of the MLP used, are⁵:

- Card: a 51-dimensional, 2-class data set based on credit approval decision [Quinlan, 1987], with 690 patterns; an MLP with 10 hidden units;
- Gene: a 120-dimensional data set with two classes, based on the detection of splice junctions in DNA sequences [Noordewier et al., 1991], with 3175 patterns; an MLP with 10 hidden units;
- Satellite: a 36-dimensional, 6-class data set with 6435 examples of feature vectors extracted from satellite imagery; an MLP with 20 hidden units.

These three sets were chosen as they have *relatively large number of features*, somewhat large number of data points, and have been studied by several researchers. Also note that the Proben1 benchmarks are particular training, validation and test splits of the UCI data sets which are available from URL <http://www.ics.uci.edu/~mlearn/MLRepository.html>. The results presented in this article are based on the first training, validation and test partition discussed in [Prechelt, 1994], where half the data is used for training, and a quarter each for validation and testing purposes.

We investigate two situations: one where the original features were randomly and disjointly partitioned among the different segments, and the second where there is some overlap among features in different segments. The exact segment count and number of features within each segment is specified in Table 5.

For each data set, we present the original number of features, the number of new features sets that result when the feature set is segmented (for Gene we only have two new sets, because the low dimensionality prevents any further segmentation), and the resulting number of features in each segment with and without overlap among the features.

A classifier trains on data from one segment, and different classifiers operate on different segments. When the number of combiners were higher than the number of segments ($N = 8$) more than on classifier (starting from a different initialization) was trained on the same features.

Tables 6-7 present the results (with the best result for each case in bold font). The misclassification percentage for individual classifiers are reported in the first column. For the trimmed mean combiner, we also provide N_1 and N_2 , the upper and lower cutting points in the ordered

⁵The number of hidden units was determined experimentally.

Table 5: Number of features in Proben1/UCI data sets

Data	Number of Original Features	Number of Segments	Features per Segment	
			no Overlap	Overlap
Card	51	4	13-13-13-12	18-18-18-18
Gene	120	4	30-30-30-30	40-40-40-40
Sat	36	4	9-9-9-9	15-15-15-15

average used in Equation 32, obtained through the validation set.

In this case, for two of the three data sets (Gene and Card), there are striking gains due to using order statistics combiners. One cause for these gains is the high variability in performance among the component classifiers. In such cases, a small number of poor classifiers can corrupt the average combiner. By their very nature, though, combiners based on order statistics are immune to this type of corruption. The ave combiner performs well on the Sat data sets where the performance among the individual classifiers is much more homogeneous. In this case, the ave results are only marginally worse than those for the trimmed mean.

Table 6: Segmented Features with overlap (% misclassified $\pm \sigma/\sqrt{n}$).

Data	N	Ave	Max	Min	Spread	Trim (N_1 - N_2)
Card	4	12.21 \pm .00	10.58 \pm .06	10.61 \pm .06	10.58 \pm .00	12.21 \pm .00 (3-4)
30.30 \pm 2.62	8	12.21 \pm .00	10.47 \pm .00	10.61 \pm .06	10.47 \pm .00	10.47 \pm .00 (7-8)
Gene	4	18.52 \pm .10	14.02 \pm .13	20.23 \pm .31	14.72 \pm .15	16.86 \pm .15 (3-4)
34.80 \pm 4.01	8	18.06 \pm .06	13.13 \pm .06	17.59 \pm .17	13.69 \pm .11	13.39 \pm .08 (7-8)
Sat	4	14.16 \pm .08	14.73 \pm .18	14.64 \pm .16	14.24 \pm .12	14.00 \pm .07 (3-4)
16.40 \pm 0.56	8	14.21 \pm .05	15.27 \pm .15	15.07 \pm .15	14.49 \pm .11	14.01 \pm .04 (3-5)

6.2 Variability through Early Stopping

For the second set of experiments we use two classes of acoustic underwater sonar signals ⁶. From the original sonar signals of four different underwater objects (porpoise sound, cracking ice and two different whale sounds), two feature sets are extracted [Ghosh et al., 1992a]:

- WOC: a 25-dimensional feature set, consisting of Gabor wavelet coefficients, temporal descriptors and spectral measurements; and,

⁶Results on 6 Proben/UCI datasets were reported in [Tumer and Ghosh, 1999] and hence are not repeated here

Table 7: Segmented Features without overlap ($\%$ misclassified $\pm \sigma/\sqrt{n}$).

Data	N	Ave	Max	Min	Spread	Trim (N_1-N_2)
Card	4	12.21 \pm .00	10.49 \pm .03	10.49 \pm .03	10.49 \pm .03	12.21 \pm .00 (3-4)
30.90 \pm 2.66	8	12.21 \pm .00	10.78 \pm .07	10.78 \pm .07	10.52 \pm .04	11.05 \pm .00 (7-8)
Gene	4	24.35 \pm .13	15.82 \pm .15	19.15 \pm .22	14.09 \pm .11	23.11 \pm .15 (3-4)
36.87 \pm 3.01	8	23.33 \pm .19	14.99 \pm .14	16.78 \pm .24	13.23 \pm .15	15.03 \pm .17 (7-8)
Sat lap	4	14.39 \pm .09	15.66 \pm .15	15.46 \pm .11	15.11 \pm .11	14.22 \pm .07 (2-3)
17.13 \pm 0.47	8	14.37 \pm .05	15.93 \pm .06	15.53 \pm .06	15.18 \pm .10	14.04 \pm .05 (3-5)

- RDO: a 24-dimensional feature set, consisting of reflection coefficients based on both short and long time windows, and temporal descriptors.

For both feature sets, an MLP with 50 hidden units was used. These data sets are available at URL <http://www.lans.ece.utexas.edu>. Further details about this 4-class problem can be found in [Ghosh et al., 1992a, Tumer and Ghosh, 1996b].

Table 8: Combining Results in the Presence of High Variability in Individual Classifier Performance for the Sonar Data ($\%$ misclassified $\pm \sigma/\sqrt{n}$).

Data	N	Ave	Max	Min	Spread	Trim (N_1-N_2)
RDO	4	11.57 \pm .11	11.94 \pm .12	11.52 \pm .20	11.04 \pm .09	11.34 \pm .14 (3-4)
13.32 \pm 0.83	8	11.64 \pm .09	11.47 \pm .11	11.29 \pm .13	11.51 \pm .09	12.30 \pm .08 (4-5)
WOC	4	8.80 \pm .09	7.84 \pm .10	9.31 \pm .12	8.54 \pm .06	8.43 \pm .13 (3-4)
12.07 \pm 1.12	8	8.82 \pm .08	7.68 \pm .12	8.91 \pm .06	8.24 \pm .11	7.81 \pm .08 (7-8)

Table 8 presents the combining results for the underwater acoustic data set when the individual classifier performance is highly variable. The results of Table 8 as well as those given in [Tumer and Ghosh, 1999] indicate that when the individual classifier performance is highly variable, order statistics-based combiners (particularly the *spread* combiner) typically provide better classification results than other simple combiners. This performance improvement is obtained without sacrificing the simplicity of the combiner. One important thing to note, however, is that in all cases studied, the order statistics based combiners performed *at least as well* as the simple combiner, implying that no risk is taken by using this method.

A close inspection of these results reveals that using either the *max* or *min* combiner can provide better classification rates than *ave*, but it is difficult to determine which of the two will

be more successful given a data set. A validation set may be used to select one over the other, but in that case, potentially precious training data is used solely for determining which combiner to use. The use of the *spread* combiner removes this dilemma by consistently providing results that are comparable to, or better than, the best of the *max-min* duo.

7 Concluding Remarks

In this article we present and analyze combiners based on order statistics. These combiners blend the simplicity of averaging with the generality of meta-learners. They are particularly effective if there are significant variations among component classifiers in at least some parts of the joint input-output space. Variations can arise when the individual training sets *cannot* be considered as random samples from a common universal data set. Examples of such cases include real-time data acquisition and classification from geographically distributed sources or data mining problems with large and possibly heterogeneous databases, where random subsampling is computationally expensive and practical methods lead to non-random subsamples [Bradley and Fayyad, 1998]. The robustness of order statistics combiners is also helpful when certain individual classifiers experience catastrophic failures (e.g., due to faulty sensors).

The analytical framework provided in this paper quantifies the reductions in error achieved when an order statistics based ensemble is used. It also shows that the two methods for linear combination of order statistics introduced in this paper provide more reliable estimates of the true posteriors than any of the individual order statistic combiners.

The experimental results of Section 5 indicate that when there is significant variability among the classifiers, the order statistics-based combiners substantially outperform simple combiners. Our previous results also showed that in the absence of such variability these combiners perform no worse. Thus the family of order statistic combiners are applicable over a wide range of situations. They are able to extract an appropriate amount of information from the individual classifier outputs without requiring tuning additional parameters as in meta-learners, and without being substantially affected by outliers.

A future endeavor, which will be helpful for this work as well as for the study of collective data mining on very large datasets in general, is to obtain a suite of public domain datasets which are intrinsically partitioned into segments with varying quality. Though such situations sometimes occur in practice (for example in oil logging data [Chakravarthy, 1997] and mortgage scoring [Merz, 1998]; both data sets proprietary), they are not represented in the standard, venerable databases such as UCI, ELENA and Statlog typically used by the academic community. Perhaps the recent Cross-Industry Standard Process for Data Mining (CRISP-DM) initiative

will provide a satisfactory solution to this problem in the near future.

Acknowledgments: This research was supported in part by ARO contracts DAAG55-98-1-0230 and DAAD19-99-1-0012, and NSF grant ECS-9900353.

References

- [Aggarwal et al., 1996] Aggarwal, J. K., Ghosh, J., Nair, D., and Taha, I. (1996). A comparative study of three paradigms for object recognition - bayesian statistics, neural networks and expert systems. In Boyer, K. and Ahuja, N., editors, *Advances In Image Understanding: A Festschrift for Azriel Rosenfeld*, pages 241–262. IEEE Computer Society Press.
- [Arnold et al., 1992] Arnold, B., Balakrishnan, N., and Nagaraja, H. (1992). *A First Course in Order Statistics*. Wiley, New York.
- [Barnett, 1981] Barnett, J. (1981). Computational methods for a mathematical theory of evidence. In *Proceedings of the Seventh International Joint Conference on Artificial Intelligence*, pages 868–875.
- [Bartlett and Shawe-Taylor, 1998] Bartlett, P. and Shawe-Taylor, J. (1998). Generalization performance of support vector machines and other pattern classifiers. In B. Scholkopf, C. J. C. B. and Smola, A., editors, *Advances in Kernel Methods - Support Vector Learning*. MIT Press, Cambridge, USA.
- [Bradley and Fayyad, 1998] Bradley, P. and Fayyad, U. M. (1998). Refining initial points for K-means clustering. In *Proceedings of the International Conference on Machine Learning (ICML-98)*, pages 91–99.
- [Breiman, 1994] Breiman, L. (1994). Bagging predictors. Technical Report 421, Department of Statistics, University of California, Berkeley.
- [Breiman, 1999] Breiman, L. (1999). Combining predictors. In Sharkey, A., editor, *Combining Artificial Neural Nets*, pages 31–50. Springer-Verlag.
- [Chakravorthy, 1997] Chakravorthy, S. (1997). Private communication. Western Atlas, Houston.
- [Chan and Stolfo, 1997] Chan, P. and Stolfo, S. (1997). On the accuracy of meta-learning for scalable data mining. *Journal of Intelligent Integration of Information*. (to appear).
- [Chan et al., 1996] Chan, P., Stolfo, S., and Wolpert (Organizers), D. (1996). Integrating multiple learned models. *Workshop with AAAI'96*.
- [Dasarathy, 1994] Dasarathy, B. (1994). *Decision Fusion*. IEEE CS Press, Los Alamitos, CA.

- [David, 1970] David, H. A. (1970). *Order Statistics*. Wiley, New York.
- [Drucker et al., 1994] Drucker, H., Cortes, C., Jackel, L. D., LeCun, Y., and Vapnik, V. (1994). Boosting and other ensemble methods. *Neural Computation*, 6(6):1289–1301.
- [DuMouchel et al., 1999] DuMouchel et al., W. (1999). Squashing flat files flatter. In Chaudhuri, S. and Madigan, D., editors, *Proceedings of the Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD-99)*, pages 6–15. ACM Press, New York.
- [Friedman, 1997] Friedman, J. H. (1997). On bias, variance, 0/1-loss, and the curse-of-dimensionality. (from: <http://www-stat.stanford.edu/~jhf/#reports>).
- [Geman et al., 1992] Geman, S., Bienenstock, E., and Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1):1–58.
- [Ghosh et al., 1992a] Ghosh, J., Deuser, L., and Beck, S. (1992a). A neural network based hybrid system for detection, characterization and classification of short-duration oceanic signals. *IEEE Journal of Ocean Engineering*, 17(4):351–363.
- [Ghosh et al., 1992b] Ghosh, J., Tumer, K., Beck, S., and Deuser, L. (1992b). Integration of local and global neural classifiers for passive sonar signals. In *Proceedings of the International Simulation Technology Conferenc*, pages 539–545, Houston, TX.
- [Granger, 1989] Granger, C. W. J. (1989). Combining forecasts—twenty years later. *Journal of Forecasting*, 8(3):167–173.
- [Hashem and Schmeiser, 1993] Hashem, S. and Schmeiser, B. (1993). Approximating a function and its derivatives using MSE-optimal linear combinations of trained feedforward neural networks. In *Proceedings of the Joint Conference on Neural Networks*, volume 87, pages I:617–620, New Jersey.
- [Holmstrom et al., 1997] Holmstrom, L., Koistinen, P., Laaksonen, J., and Oja, E. (1997). Neural and statistical classifiers - taxonomy and two case studies. *IEEE Transactions on Neural Networks*, 8:5–17.
- [Jordan and Jacobs, 1994] Jordan, M. I. and Jacobs, R. A. (1994). Hierarchical mixture of experts and the EM algorithm. *Neural Computation*, 6:181–214.
- [Kanal, 1974] Kanal, L. (1974). Patterns in pattern recognition. *IEEE Transactions on Information Theory*, 20:697–722.
- [Kargupta and Park, 2000] Kargupta, H. and Park, B.-H. (2000). Collective data mining: A new perspective toward distributed data mining. In Kargupta, H. and Chan, P., editors, *Advances in Distributed and Parallel Knowledge Discovery*. AAAI/MIT Press, Cambridge, MA.

- [Krogh and Vedelsby, 1995] Krogh, A. and Vedelsby, J. (1995). Neural network ensembles, cross validation and active learning. In Tesauro, G., Touretzky, D. S., and Leen, T. K., editors, *Advances in Neural Information Processing Systems-7*, pages 231–238. M.I.T. Press.
- [Merz and Pazzani, 1997] Merz, C. and Pazzani, M. (1997). Combining neural network regression estimates with regularized linear weights. In Mozer, M. C., Jordan, M. I., and Petsche, T., editors, *Advances in Neural Information Processing Systems-9*, pages 564–570. M.I.T. Press.
- [Merz, 1998] Merz, C. J. (1998). Private communication. HNC Software, San Diego.
- [Minsky, 1991] Minsky, M. (1991). Logical versus analogical or symbolic versus connectionist or neat versus scruffy. *AI Magazine*, 12(2):34–51.
- [Mitchell, 1997] Mitchell, T. (1997). *Machine Learning*. Morgan Kaufmann, San Mateo, CA.
- [Narendra et al., 1995] Narendra, K., Balakrishnan, J., and Ciliz, K. (1995). Adaptation and learning using multiple models, switching and tuning. *IEEE Control Systems Magazine*, pages 37–51.
- [Noordewier et al., 1991] Noordewier, M. O., Towell, G. G., and Shavlik, J. W. (1991). Training knowledge-based neural networks to recognize genes in DNA sequences. In Lippmann, R., Moody, J., and Touretzky, D., editors, *Advances in Neural Information Processing Systems-3*, pages 530–536. Morgan Kaufmann.
- [Oates and Jensen, 1998] Oates, T. and Jensen, D. (1998). Large datasets lead to overly complex models: An explanation and a solution. In *Proc. KDD-98*, pages 294–298.
- [Perrone and Cooper, 1993] Perrone, M. and Cooper, L. N. (1993). When networks disagree: Ensemble methods for hybrid neural networks. In Mammone, R. J., editor, *Neural Networks for Speech and Image Processing*, chapter 10. Chapman-Hall.
- [Perrone, 1993] Perrone, M. P. (1993). *Improving Regression Estimation: Averaging Methods for Variance Reduction with Extensions to General Convex Measure Optimization*. PhD thesis, Brown University.
- [Prechelt, 1994] Prechelt, L. (1994). PROBEN1 — A set of benchmarks and benchmarking rules for neural network training algorithms. Technical Report 21/94, Fakultät für Informatik, Universität Karlsruhe, D-76128 Karlsruhe, Germany. Anonymous FTP: /pub/papers/tech-reports/1994/1994-21.ps.Z on ftp.ira.uka.de.
- [Prodromidis et al., 2000] Prodromidis, A., Chan, P., and Stolfo, S. (2000). Meta-learning in distributed data mining systems: Issues and approaches. In Kargupta, H. and Chan, P.,

- editors, *Advances in Distributed and Parallel Knowledge Discovery*. AAAI/MIT Press, Cambridge, MA.
- [Provost, 2000] Provost, F. (2000). Distributed data mining: Scaling up and beyond. In Kargupta, H. and Chan, P., editors, *Advances in Distributed and Parallel Knowledge Discovery*. AAAI/MIT Press, Cambridge, MA.
- [Provost and Kolluri, 1999] Provost, F. and Kolluri, V. (1999). A survey of methods for scaling up inductive learning algorithms. *Data Mining and Knowledge Discovery Journal*, 3(2):131 – 169.
- [Quinlan, 1987] Quinlan, J. (1987). Simplifying decision trees. *International Journal of Man-Machine Studies*, 27:221–234.
- [Ramamurti and Ghosh, 1999] Ramamurti, V. and Ghosh, J. (1999). Structurally adaptive modular networks for non-stationary environments. *IEEE Transactions on Neural Networks*, 10(1):152–60.
- [Richard and Lippmann, 1991] Richard, M. and Lippmann, R. (1991). Neural network classifiers estimate Bayesian a posteriori probabilities. *Neural Computation*, 3(4):461–483.
- [Ruck et al., 1990] Ruck, D. W., Rogers, S. K., Kabrisky, M. E., Oxley, M. E., and Suter, B. W. (1990). The multilayer Perceptron as an approximation to a Bayes optimal discriminant function. *IEEE Transactions on Neural Networks*, 1(4):296–298.
- [Sarhan and Greenberg, 1956] Sarhan, A. E. and Greenberg, B. G. (1956). Estimation of location and scale parameters by order statistics from singly and doubly censored samples. *Annals of Mathematical Statistics Science*, 27:427–451.
- [Schapire et al., 1997] Schapire, R., Freund, Y., Bartlett, P., and W.S., L. (1997). Boosting the margin: A new explanation for the effectiveness of voting methods. In *Proceedings of the Fourteenth International Conference on Machine Learning*. Morgan Kaufmann.
- [Selfridge, 1958] Selfridge, O. G. (1958). Pandemonium: a paradigm for learning. *Proc. of Symp. held at the National Physical Lab.*, pages 513–526.
- [Sharkey, 1996] Sharkey, A. J. J. (1996). On combining artificial neural nets. *Connection Science, Special Issue on Combining Artificial Neural Networks: Ensemble Approaches*, 8(3 & 4):299–314.
- [Shepherd, 1979] Shepherd, G. (1979). *The Synaptic Organization of the Brain*. Oxford Univ. Press, 2 edition.

- [Taha and Ghosh, 1997] Taha, I. and Ghosh, J. (1997). Hybrid intelligent architecture and its application to water reservoir control. *International Journal of Smart Engineering Systems*, 1:59–75.
- [Tumer and Ghosh, 1996a] Tumer, K. and Ghosh, J. (1996a). Analysis of decision boundaries in linearly combined neural classifiers. *Pattern Recognition*, 29(2):341–348.
- [Tumer and Ghosh, 1996b] Tumer, K. and Ghosh, J. (1996b). Error correlation and error reduction in ensemble classifiers. *Connection Science, Special Issue on Combining Artificial Neural Networks: Ensemble Approaches*, 8(3 & 4):385–404.
- [Tumer and Ghosh, 1999] Tumer, K. and Ghosh, J. (1999). Linear and order statistics combiners for pattern classification. In Sharkey, A. J. C., editor, *Combining Artificial Neural Nets: Ensemble and Modular Multi-Net Systems*, pages 127–162. Springer-Verlag, London.
- [Wolpert and Tumer, 1999] Wolpert, D. and Tumer, K. (1999). A survey of collective intelligence. In Bradshaw, J. M., editor, *Handbook of Agent Technology*. AAAI Press/MIT Press, 1999, Cambridge, MA.
- [Wolpert, 1992] Wolpert, D. H. (1992). Stacked generalization. *Neural Networks*, 5:241–259.